



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2021

Eliciting beliefs as distributions in online surveys

Leemann, Lucas ; Stoetzer, Lukas F ; Traunmueller, Richard

Abstract: Citizens' beliefs about uncertain events are fundamental variables in many areas of political science. While beliefs are often conceptualized in the form of distributions, obtaining reliable measures in terms of full probability densities is a difficult task. In this letter, we ask if there is an effective way of eliciting beliefs as distributions in the context of online surveys. Relying on experimental evidence, we evaluate the performance of five different elicitation methods designed to capture citizens' uncertain expectations. Our results suggest that an elicitation method originally proposed by Manski (2009) performs well. It measures average citizens' subjective belief distributions reliably and is easily implemented in the context of regular (online) surveys. We expect that a wider use of this method will lead to considerable improvements in the study of citizens' expectations and beliefs.

DOI: <https://doi.org/10.1017/pan.2020.42>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-199123>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) License.

Originally published at:

Leemann, Lucas; Stoetzer, Lukas F; Traunmueller, Richard (2021). Eliciting beliefs as distributions in online surveys. *Political Analysis*, 29(4):541-553.

DOI: <https://doi.org/10.1017/pan.2020.42>

Eliciting Beliefs as Distributions in Online Surveys^{*}

Lucas Leemann¹, Lukas F. Stoetzer², and Richard Traunmüller³

¹Department of Political Science, University of Zürich, Switzerland. Email: leemann@ipz.uzh.ch

²Department of Political Science, Humboldt University of Berlin, Germany. Email: lukas.stoetzer@hu-berlin.de

³School of Social Sciences, University of Mannheim, Germany. Email: traunmueller@uni-mannheim.de

Abstract

Citizens' beliefs about uncertain events are fundamental variables in many areas of political science. While beliefs are often conceptualized in the form of distributions, obtaining reliable measures in terms of full probability densities is a difficult task. In this letter, we ask if there is an effective way of eliciting beliefs as distributions in the context of online surveys. Relying on experimental evidence, we evaluate the performance of five different elicitation methods designed to capture citizens' uncertain expectations. Our results suggest that an elicitation method originally proposed by Manski (2009) performs well. It measures average citizens' subjective belief distributions reliably and is easily implemented in the context of regular (online) surveys. We expect that a wider use of this method will lead to considerable improvements in the study of citizens' expectations and beliefs.

Keywords: prior elicitation, online survey research, citizen beliefs, measurement.

1 Introduction

Citizens' beliefs about uncertain events are fundamental variables in many areas of political science, including work on attitudes (e.g. Zaller and Feldman 1992), cognitive biases (e.g. Gerber and Green 1999; Bartels 2002; Bullock 2009), ambivalence (e.g. Alvarez and Brehm 1997), misinformation (e.g. Berinsky 2017), or citizen forecasts (e.g. Murr 2011; Leiter, Ramirez, and Stegmaier 2018), to name just a few. While beliefs are often theoretically conceptualized in the form of distributions, obtaining reliable measures of these beliefs in terms of full probability densities is a difficult task (Savage 1971; Garthwaite, Kadane, and O'Hagan 2005; Goldstein and Rothschild 2014). Most survey questions are focused on the first moment of an underlying distribution and thus miss important information about beliefs' variance or uncertainty.

The question we ask in this letter is whether there is an effective way to elicit average citizens' belief distributions in the context of online surveys? This paper discusses five different elicitation methods designed to capture citizens' uncertain expectations. We present experimental evidence and evaluate which question format is best suited to elicit continuous beliefs as distributions from regular (i.e. non-expert) survey respondents. That is, we are interested in how well these methods capture subjective distributions when compared to a benchmark and which of these methods performs best.

Our results suggest that an elicitation method originally proposed by Manski 2009 performs well. It contains five sequential survey questions that reliably measure average citizens' subjective

Political Analysis (2020)

DOI: 10.1017/pan.xxxx.xx

Corresponding author
Lucas Leemann

Edited by
John Doe

© The Author(s) 2020. Published
by Cambridge University Press
on behalf of the Society for
Political Methodology.

^{*} Many thanks to Daniel Bischof, Tim Hicks, Patrick Kraft, Andreas Murr, Simon Munzert, Ana Petrova, and David Rothschild for their valuable comments and help. We would also like to thank the editor for his guidance and help and the two anonymous referees for their comments. We thank Lucien Baumgartner for his impeccable research assistance. This research was partly funded by the Department of Political Science at the University of Zürich.

belief distributions and that are easily implemented in the context of regular online surveys. They are also easy and quick to answer and, hence, not too cost-intensive in online surveys. We expect that a wider use of this method will lead to considerable improvements in the study of citizens' expectations and beliefs and, therefore, to important political science theories. In addition, it should also prove a useful tool to Bayesians who wish to elicit subjective prior distributions from non-experts (Gill and Walker 2005).

To illustrate the use of the method in an applied example, we elicit people's expectations about the 2020 U.S. presidential election. Eliciting citizens' beliefs is a common element in citizen forecasts (Murr 2011), for which it would be valuable to distinguish between citizens who are more certain (i.e. who have narrow belief distributions) from those who are less certain about the election outcome (i.e. who have wide belief distributions) and weight them accordingly. Hence, we ask respondents to provide their full belief distribution concerning Donald Trump's likely vote share in the November 2020 election. In section 5 we describe how the elicitation methods discussed in this letter can be applied to this practical task. Based on the Manski question format, we find that respondents expect a popular vote share of 48 % for Donald Trump with a standard deviation of 5.5%. We further find considerable differences in both expectations and uncertainties between Democrats (44 %, sd 4.6%) and Republicans (52%, sd 8.1%).

The remainder of this letter proceeds as follows. The next section discusses the elicitation process. Section 3 then presents the experimental setup and the five elicitation approaches we evaluate. Section 4 presents the results. Section 5 provides a brief illustration using Trump's vote share in the November 2020 election as an example. Section 6 concludes.

2 Eliciting Beliefs as Distributions

The elicitation of beliefs as distributions has a long tradition in statistics, psychology, and economics. In political science, Bayesians seek to elicit prior distributions from *experts* to inform their statistical models (Gill and Walker 2005; Gill and Freeman 2013). However, the process of eliciting probability distributions described in this literature usually is a time-consuming enterprise that requires careful effort even when it is used to learn about the beliefs of experts who may already be familiar with probabilities.

What makes the elicitation of beliefs so difficult is that average people are not used to expressing themselves in easily quantifiable ways. Many citizens are unlikely to be familiar with the concept of probability and not used to expressing their expectations in terms of distributions. Lengthy elicitation protocols also do not scale well to the number of respondents required for testing political science theories about citizens' expectations and are unlikely to be part of nationally representative surveys. Thus, the central challenge is how to best translate what people think into probability distributions within the confines of standard survey methodology.

Formally, an elicitation process can involve up to four steps (Garthwaite, Kadane, and O'Hagan 2005). In the *setup* step, the problem is defined and respondents are recruited and trained in the key concepts and procedures. *Elicitation* is the key step where the respondent is asked to provide information about his or her subjective belief. In the *fitting* step, this elicited information is converted into a probability distribution. The final step assesses the *adequacy* of the elicited distribution and provides an opportunity for correction. The challenge we address in this letter is how to implement these steps in the context of regular online surveys, where time and scale concerns as well as limited researcher-respondent interaction render the use of full elicitation protocols impractical.

Traditional elicitation methods come in three basic forms (Spetzler and Stael von Holstein 1975). In each of these three forms, subjects are asked questions and the answers represent points on a cumulative distribution function. In so-called P-methods, subjects are provided with fixed values referring to the quantity of interest and asked to assign *probabilities* attached to these values (e.g.

what is the probability that the value is below x ?). In V-methods, subjects are instead provided with pre-defined probabilities and asked to assign *values* to them (e.g. at what value are half of the observations below or above that value?). PV-methods are more difficult and simultaneously integrate both approaches. For instance, respondents may be asked to draw a graph of a probability distribution. In this letter, we evaluate several ways to implement these methods with online survey questions.

Given humans' difficulties with probabilities, eliciting beliefs as distributions is as much a psychological problem as it is a statistical one. Many cognitive human biases are well known: representativeness, availability, anchoring biases, the law of small numbers as well as hindsight biases (Tversky and Kahneman 1971, 1973, 1974; Kynn 2008). But it is important to distinguish those biases in beliefs from biases introduced by elicitation methods. Psychological research suggests that while people are generally capable of estimating proportions, modes, and medians, they are less proficient at assessing the means of highly skewed distributions (Peterson and Miller 1964) and often have serious misconceptions about variances (Garthwaite, Kadane, and O'Hagan 2005). People are reasonably good at quantifying their opinions as credible intervals but have the tendency to imply a greater degree of confidence than is justifiable (Wallsten and Budescu 1983; Cosmides and Tooby 1996).

3 Experimental Set-Up

In the following, we evaluate a set of elicitation question formats. For a proper evaluation of elicitation methods we need an objective benchmark against which to judge the derived beliefs. To this end, we run a number of experiments where we instill objective distributions and assess which format yields beliefs that are most consistent with these objective benchmark distributions.

Instead of working with arbitrary numbers, we rely on an example of citizens' beliefs about hypothetical election results. Note that this experimental evaluation is different from an actual elicitation process where we would not instill a prior but rather try to elicit a pre-existing belief. To illustrate the actual usage of the method to a political science audience, we provide an example of an actual elicitation process further below. In the following presentation of our experiments, we proceed along the four steps of the elicitation process described in the previous section: setup, elicitation, fitting, and adequacy check.¹

3.1 The Setup Step

We ran experiments with a total of about 3,600 participants. We relied on Amazon Mechanical Turk (MTurk), which is widely used for scientific purposes (Berinsky, Huber, and Lenz 2012; Mason and Suri 2012; Thomas and Clifford 2017). We recruited workers advertising a study on *surveys*, *opinion polls*, and *charts*. MTurk allowed us to carry out the experiments in a short time period and at a low cost. While MTurk samples may be special, they are comparable to other online samples. Mullinix *et al.* 2015 analyze treatment effects obtained from 20 experiments implemented on a population-based sample and MTurk. The results reveal considerable similarity between effects obtained from convenience and nationally representative population-based samples. Coppock 2018 replicates fifteen survey experiments and compares the estimates based on random samples to estimates based on an MTurk sample. In general, the two sets of estimates overlap. These findings may not be surprising because just like MTurk, many online survey panels actually consist of semi-professional survey takers who are experienced in completing online tasks and are perhaps younger and more educated (see e.g. Berinsky, Huber, and Lenz 2012, p.358), in addition to being paid for their participation. Since we are specifically interested in eliciting beliefs in the context of

1. An "adequacy check" gives respondents the chance to review and check their elicited belief distribution and correct themselves in case this elicited distribution does not adequately represent their belief.

online surveys, these particular respondent characteristics do not concern us much.

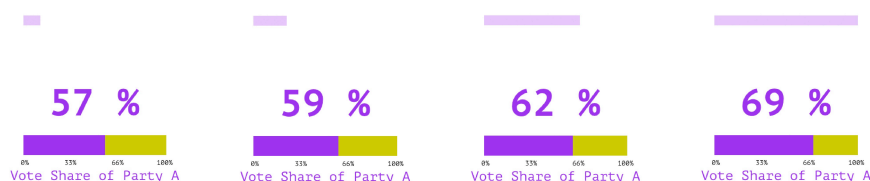


Figure 1. Four still frames from the GIF. Each still frame is shown for about half a second. The black bar at the top is a progress bar.

We presented respondents with 100 results from hypothetical local elections that we randomly drew from a pre-specified distribution. By exposing respondents to these draws, we manipulated the objective belief distribution² along two factors: symmetric vs. asymmetric and small vs. large variance. We rely on a Beta distribution in all four conditions but vary the shape parameters of the distribution. The symmetric small-variance distribution is $\mathcal{B}(60,60)$ and the respective asymmetric distribution is $\mathcal{B}(60,30)$. For the large-variance condition, we rely on $\mathcal{B}(30,30)$ for the symmetric and on $\mathcal{B}(30,15)$ for the asymmetric distribution. The hypothetical results of 100 election simulations are presented as a short GIF where each frame shows one election outcome and is displayed for about half a second. Four random draws are illustrated in Figure 1. This approach follows Goldstein and Rothschild 2014, who also rely on this form of visualization to present the distribution. The goal is to treat these distributions as the objective truth and to identify which question format elicits beliefs that are closest to the true distribution.

Each respondent is then also randomly assigned to an elicitation question format in a simple *between-subjects design* (one question format per respondent). There is balance across question types with respect to a number of socio-economic variables (see section 3 in the appendix). We also employ two questions that serve as attention checks, and each question is correctly answered by about 75% of the respondents. Here, we show results for all respondents that answered both questions correctly, which is about 60% of the original sample. The same tables based on all respondents are shown in the appendix (see section 5 in the appendix). There is no substantive difference between the two.

3.2 The Elicitation Step: Comparing Five Question Formats

The literature proposes different question formats to elicit univariate distributions (e.g. O'Hagan *et al.* 2006, chapter 5.2). Here, we compare five common question formats that elicit different elements of a distribution and pose varying levels of cognitive demand. Two main selection criteria guided our choice of formats: a) general question type and b) ease of implementation in the context of online surveys. Based on a review of the relevant literature, we found that different question formats refer to different aspects of the belief distribution. Some present fixed intervals and ask for probabilities, others directly elicit quantile values or rely on a mix of both (see the distinction of P-methods and V-methods mentioned above). While most elicitation methods are purely verbal, others make use of visualization. Thus, our goal was to include one method of each general type.

Equally important is the second goal: to evaluate only such methods that are easily implemented in online surveys because they follow a simple question format. In addition, we also take advantage of the fact that online surveys provide us with the ability to use simple visual tools. But we will not consider elicitation protocols that demand close researcher-respondent interaction (e.g. Morris, Oakley, and Crowe 2014) or rely on incentivized elicitation methods that are often used in economic

2. This works under the assumption that in our hypothetical example, respondents hold no priors over the result.

laboratory experiments (for an overview, see e.g. Schlag, Tremewan, and Van der Weele 2015).

Here, we only briefly discuss each format. We present precise question wording in the appendix (section 1).

Interval Question (Wide and Narrow). These questions ask about the *probabilities* of fixed intervals (with the two versions varying the width of the interval values). More specifically, respondents are first asked to indicate the most likely value and then to provide us with the probability that a vote outcome will be lower than 40% (45% in the narrow format) and the probability that it will be higher than 60% (55% in the narrow format).

Quantile Question. The second question format asks respondents to provide three quantile *values*: the median, the first quartile and the third quartile. This question format also provides an adequacy check. It ends by showing people their three responses (P_{25} , P_{50} , P_{75}) and asking them whether they think that a random draw is equally likely to fall into any of these intervals: $0-P_{25}$, $P_{25}-P_{50}$, $P_{50}-P_{75}$, $P_{75}-1$. Respondents can then correct themselves if they wish to do so. Thus, the fourth elicitation step is possible and respondents can assess the adequacy of their responses.

Manski Question. The third hybrid question format relies on work by Manski 2009 and asks for both, values *and* probabilities. Specifically, it first asks for three values along the distribution (the most likely value as well as the expected lower and upper bounds) and then asks respondents to give provide probabilities for their elicited lower and upper bounds.

Bins and Balls. The last question is the latest addition to elicitation methods and takes advantage of the fact that a large number of surveys are being carried out online and, hence, allow for completely new question formats. Bins and Balls follows a proposal by Goldstein and Rothschild 2014 and is a visual tool for specifying a distribution where respondents have to place 100 balls into bins of a specific range (see Figure 2). Balls are placed in a bin by the respondent's clicking on the + and – symbols. Since respondents are able to directly see the implied distribution, this is akin to an implicit adequacy check.

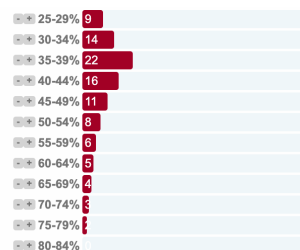


Figure 2. Screenshot of a Balls and Bins question. Illustration after hypothetical respondent has allocated all 100 balls.

All five question formats differ in their complexity for respondents but also in how easily they can be implemented. Some of these questions lend themselves to adding an adequacy check at the end, others do not. Table 1 allows us to compare the different formats. The number of questions is ill-defined for the Bins and Balls format as it is one question but requires respondents to provide 100 inputs to distribute the virtual balls.

The Interval Question and the Quantile Question are particularly demanding as they require an understanding of quantiles. The Manski method is similar but can be expected to be less demanding since it translates the task into easier terms (means, maximums, and minimums) well. Finally, the Bins and Balls Question requires the least of respondents but its implementation is the most

Table 1. Question Formats

Question format	Most difficult concept R's need to know?	How many questions?	Adequacy check?
Interval (Wide and Narrow)	Quantiles	3	×
Quantile	Median	4	✓
Manski	Percentages	5	×
Bins & Balls	Percentages	1/100	✓

demanding for researchers. The question formats further differ on whether they allow for an adequacy check. In the Quantile question, for example, respondents can incorrectly place the upper quartile below the median. This can signal a wrong understanding of the question. In the next section, we investigate the accuracy of the elicited beliefs and discuss the experimental results.

3.3 The Fitting Step

To estimate a respondent's belief, we assume a flexible parametric distribution for his or her beliefs and estimate the parameters of the distribution such that it closely mimics the observed indicators for the different question formats. Because the sampling space of our experiment is bound between 0 and 1, we employ a Beta distribution as our parametric assumption. The Beta distribution has two shape parameters: α and β . We provide the derivation of the interval question format as an example here. We present the derived likelihood functions for the other formats in the appendix (see section 2).

We observe three values for the interval question. Respondents report the mean value of their beliefs and the probabilities of observing a value below and above a certain threshold. We denote the mean with y_i and the two ($k \in (1, 2)$) probabilities with p_{i1} and p_{i2} . The interval values depend on the question format and are denoted with $c = [c_1, c_2]$, where in the wide version $c = [40\%, 60\%]$ and in the narrow version $c = [45\%, 55\%]$. We assume that the values are measured with normal measurement error.³

$$y_i \sim \mathcal{N}(\mu_y, \sigma_y^2) \quad (1)$$

$$p_{i1} \sim \mathcal{N}(\mu_{p_1}, \sigma_p^2) \quad (2)$$

$$p_{i2} \sim \mathcal{N}(\mu_{p_2}, \sigma_p^2) \quad (3)$$

The expectations μ_y are calculated from the assumed parametric belief distribution. Here, we use the same distribution as in the data-generating process - a beta distribution. The beta distribution is relatively flexible and well-suited for our example with vote shares constrained on the unit interval. The expectation for the mean from the beta is given by the two shape parameters α and β :

$$\mu_y = \frac{\alpha}{\alpha + \beta} \quad (4)$$

The expected probabilities are given by the CDF of the beta distribution, which we denote with $Q(\cdot, \alpha, \beta)$.

3. Including a measurement model extends existing approaches that only minimize the squared error between observed and theoretical expected value (e.g. Morris, Oakley, and Crowe 2014). This can open the modelling framework up for a set of extensions, e.g. correlated and heteroscedastic errors, hierarchical structures, and Bayesian estimation.

$$\mu_{p_1} = Q(c_1, \alpha, \beta) \quad (5)$$

$$\mu_{p_2} = 1 - Q(c_2, \alpha, \beta) \quad (6)$$

With this model, we can define the Likelihood for the observed data. As we assume that all responses are identically and independently normal distributed, the likelihood is the product of three normal distributed measurements y_i , p_{i1} and p_{i2} for each of the respondents.⁴ To obtain Maximum Likelihood estimates of the parameters α , β , σ_p , σ_y , the log-likelihood function is maximized using R's `optim` function. The estimates yield an estimate of the average beliefs for a specific condition. In the experiments, we can then identify the question format that will yield average belief estimates closest to the true values.

3.4 The Adequacy Check Step

Assessing the adequacy of the elicited distribution by giving respondents the chance to review and correct their belief distributions is difficult, because the fitting is done 'outside' of the survey software and only after the answers have been collected. But for some formats, it is still possible to provide the opportunity for correction using question filters based either on respondents' answers or on visual question formats. The Quantile Question, for instance, presents respondents with the quartiles they provided and asks if election results are equally likely to fall within each of them.⁵ The Bins and Balls format asks respondents to 'draw' their distribution and thus provides immediate feedback.

4 Results

To evaluate the different elicitation methods, we now compare the elicited beliefs to the benchmark of true objective distributions. Each column in [Figure 3](#) stands for a combination of conditions (small/large variance and symmetric/asymmetric distribution). While we look at both symmetric and asymmetric true distributions, the asymmetric scenarios are likely to be more relevant in practice. This is because the only symmetric Beta distributions are those distributions where the two shape parameters are exactly equal to each other. The five rows contain the different elicitation methods. We focus on the *average* elicited belief across all respondents and present the same figure with each *individual* belief distribution in the Appendix (see Figure X).⁶

We find that most question formats are unbiased when the true distribution is symmetric, i.e. they are able to provide the correct first moment. With asymmetric distributions, there is some bias towards .5 but its extent varies across question formats. It is especially evident for the Bins and Balls format. Looking at the second moment, we find that the two Interval questions tend to provide beliefs that are too wide in both, the symmetric and asymmetric scenarios. Thus, after simply eyeballing the plots, it seems that overall the Manski question and the Quantile question come closest to the true distributions.

To evaluate the question formats more formally, we turn to the results in [Table 2](#) where we illustrate for each combination of experimental factors: the implied parameters of the elicited priors, the Kullback-Leibler divergence, the number of observations, and the p-value of a likelihood-ratio test on whether the estimated parameters differ from the true values of the parameter. The smaller the KL divergence, the closer the elicited prior is to the true distribution. We also present a figure with the sum of the KL divergence over all four experimental conditions (see [Figure 4](#)). Here, we

4. We provide a more detailed description of the Likelihood function in section A2 in the appendix.

5. In Appendix (section 4), we compare the results of the Quantile question with and without the adequacy check and find that the adequacy check can considerably improve the result.

6. A full replication package is available, see Leemann, Traunmueller, and Stoetzer 2020.

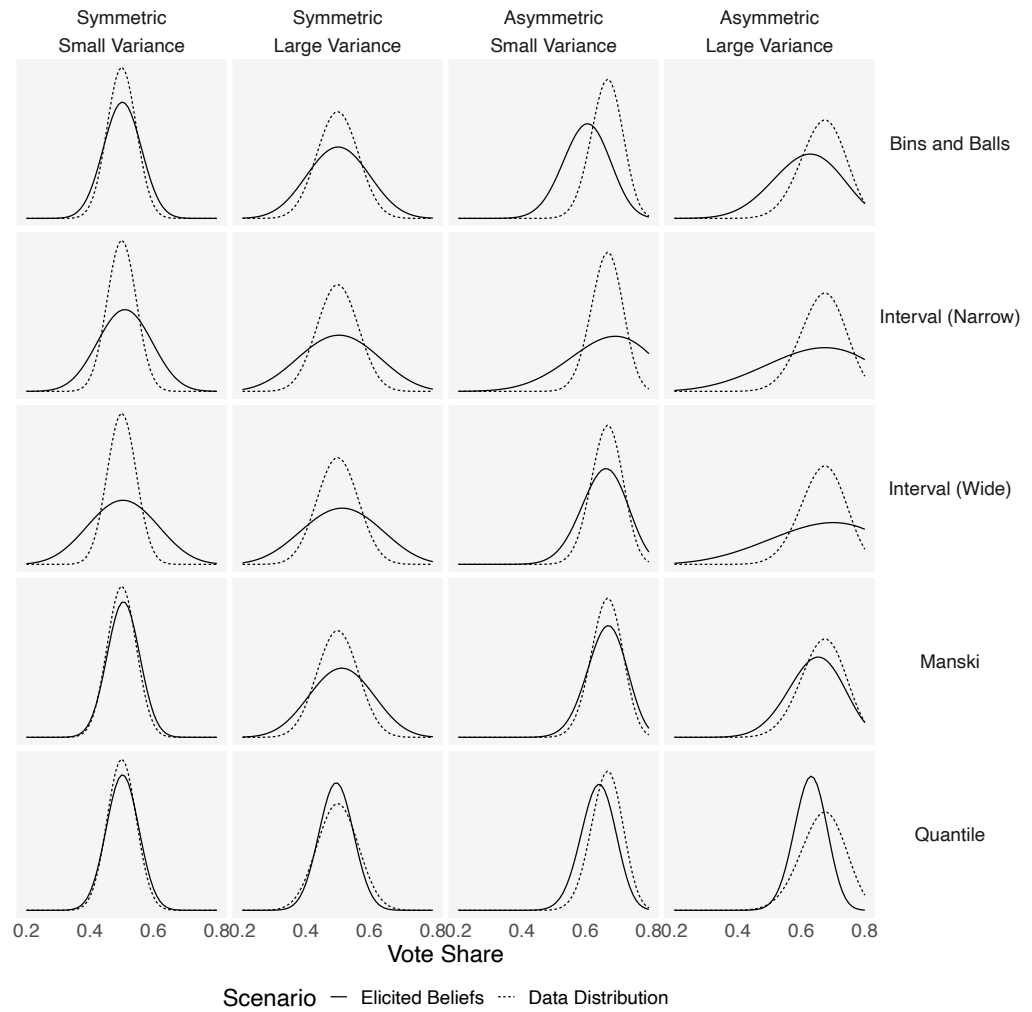


Figure 3. Comparison of Question Formats. The dotted line indicates the true distribution and the black solid line shows the average of the elicited distributions.

again only show the results when averaging across all respondents. The appendix contains the results for individual respondents (see Figure 3).

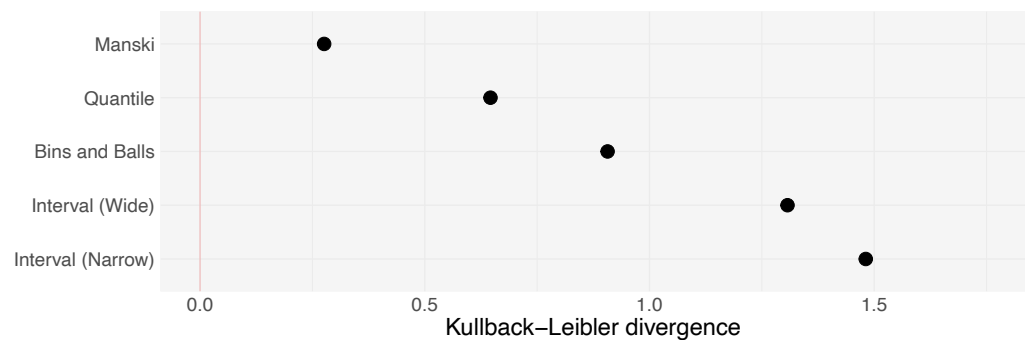


Figure 4. Summed Kullback-Leibler Divergence. The Manski question performs best across the five experimental scenarios.

method	alpha	beta	KL	lr	N
Quantile	42.35	43.10	0.04	0.44	65
Bins and Balls	13.59	13.51	0.13	0.00	60
Manski	13.01	12.41	0.15	0.00	62
Interval (Narrow)	8.54	8.41	0.28	0.00	61
Interval (Wide)	8.68	8.27	0.29	0.00	69

(a) Symmetric, Large Variance

method	alpha	beta	KL	lr	N
Quantile	48.41	47.90	0.01	0.69	119
Manski	48.75	47.72	0.02	0.39	112
Bins and Balls	35.67	35.38	0.06	0.00	107
Interval (Narrow)	18.02	17.35	0.27	0.00	126
Interval (Wide)	11.03	10.88	0.45	0.00	115

(c) Symmetric, Small Variance

method	alpha	beta	KL	lr	N
Manski	19.97	11.07	0.06	0.00	65
Bins and Balls	12.75	7.96	0.23	0.00	48
Interval (Narrow)	5.95	3.39	0.41	0.00	68
Interval (Wide)	5.39	2.88	0.46	0.00	65
Quantile	55.23	32.69	0.47	0.00	65

(b) Asymmetric, Large Variance

method	alpha	beta	KL	lr	N
Manski	38.63	19.36	0.04	0.18	112
Interval (Wide)	28.39	14.79	0.11	0.00	120
Quantile	49.08	27.64	0.13	0.00	118
Bins and Balls	27.21	17.98	0.50	0.00	121
Interval (Narrow)	9.38	4.69	0.52	0.00	133

(d) Asymmetric, Small Variance

Table 2. Experimental Comparison of Five Elicitation Methods Across Four Scenarios. Implied parameters of elicited priors, Kullback-Leibler divergence, p-value of a likelihood-ratio test and number of observations shown.

If we take the sum of all four experimental settings, the Manski question scores the smallest value for the Kullback-Leibler divergence, $KL = 0.28$. It is followed by the Quantile question ($KL = 0.65$) and Bins and Balls ($KL = 0.91$). The two Interval questions perform worst ($KL = 1.31$ for the wide and $KL = 1.48$ for the narrow interval.) As mentioned above, in practice asymmetric scenarios are more frequent and the Manski question format also beats all other alternatives for this case.

Based on these experiments we conclude that the Manski question format outperforms the other elicitation methods. In principle, it seems intuitive that eliciting more points along a distribution would also result in a better measurement of respondents' beliefs.⁷ However, we would be mistaken to equate the number of questions with an elicitation's methods performance. The Quantile question, for instance, asks for three quantities and adds an adequacy check. Without this adequacy check (an option that we test, see Table A5 in the appendix), it would ask just as many questions as the two Interval questions - yet the performance across these formats clearly differs. Without the adequacy check, the Quantile question outperforms the two Interval questions in the symmetric scenario with large variance (KL-distance of .07 vs. .28 and .29) but has more problems than the two Interval questions in the asymmetric scenario with large variance (KL-distance of .98 vs. .41 and .46). One possibility why the Quantile question fares better in the symmetric case than in the asymmetric case is that in the symmetric case respondents can rely on equal distance of the first and third quantile to the median as an informal adequacy check. This possibility does not exist when the belief is asymmetric.

In sum, the Manski format provides a fairly effective approach to prior elicitation that is straightforward to implement because it only requires five questions to measure respondents' beliefs. In addition, the Manski question takes marginally less time (median completion time was 170 seconds) than the Quantile question (200s) and Bins and Balls (199s), but more time than the two Interval questions (149 and 157s, respectively). These differences are not statistically significant (more detail is provided in section A8 in the appendix).

A valid question is whether our results are sensitive to the composition of the sample used.

7. Note that this is different from the classical notion of reducing measurement error by have multiple measures of the same underlying quantity.

For example, it is unclear whether respondents on MTurk pay more or less attention than 'normal' respondents. While some argue that MTurkers are less attentive and try to complete tasks as quickly as possible, others argue that workers are more attentive because they are paid for these tasks. Several studies have looked into the properties of MTurk samples and found them to perform equally well to other online samples (Mullinix *et al.* 2015; Coppock 2018). We provide analyses that probe into the effects of respondents' attention in the Appendix (section A5). Comparing attentive respondents (i.e. those that passed the attention checks⁸) to all respondents, we find that attentive respondents are slightly closer to the objective distributions than the complete sample. More importantly however, the relative performance of the five elicitation methods is not affected by respondents' level of attentiveness. We find similar results for the distinction between sophisticated and unsophisticated respondents (as proxied by political interest, see section A6 in the Appendix).

On a final note, we only find limited evidence for any systematic biases in respondents' beliefs. In particular, we only observe over-confidence (i.e. respondents' tendency to be more certain than the objective data would warrant and, therefore, assign variances that are too narrow) in the case of the Quantile question. In the symmetric scenario with large variance, for instance, the true standard deviation is $\sigma = .064$, but the average elicited distribution yields a standard deviation of only $\sigma = .051$.⁹ For all other formats, respondents actually express beliefs that are *less* certain than the objective benchmark would demand (i.e. the elicited distributions are too wide).

Before concluding this letter, we provide an illustrating application where we use these different techniques to elicit people's subjective beliefs about a future outcome.

5 Application: What Vote Share Will Trump Receive in November 2020?

In the applied setting of an actual elicitation process, researchers would of course not instill an objective distribution. Instead they would seek to elicit the beliefs that respondents already hold about a subject matter. To illustrate the relative performance of the five elicitation methods in a more realistic setting, this section provides an example where we ask respondents to indicate their beliefs about the upcoming presidential election. This closely mimics an actual elicitation exercise.

We report the results from an online survey carried out on MTurk with 500 participants. Each respondent was asked what their belief was of the popular vote share that Donald Trump will garner in November 2020. As with the main experiments presented above, we again only offer respondents one randomly assigned question format in a simple between-subjects design. For each question format, we estimate the underlying belief distribution of the full sample and then, because we expect clear partisan differences, separately for Democrats and Republicans.

There are two main results that can be gleaned from Figure 5 (we provide more detail on the estimated beliefs in Table A10 in the online appendix). First, which question format is used for eliciting respondents' beliefs clearly matters. The average expected 2020 popular vote share for Donald Trump differs from one format to the other. In addition, the elicitation methods differ vastly in the belief variances they produce (see section A7 in the appendix for more details). In the experiments, the Manski question format performed best in retrieving objective belief distributions. When eliciting pre-existing subjective beliefs we do not know the true distribution and hence cannot assess each method's precision. What we can assess, is how clearly the signal is measured, i.e. which formats provide plausible results with low variance. In the 2020 Trump election example, we find that the ordering in performance is similar to the ordering found in the experiments. Given

8. We use two screening questions to detect inattentive respondents (Berinsky, Margolis, and Sances 2014). One asks respondents to recall the founding organization mentioned in the introductory text and the second question asks them to recall the colors of the plots in the GIF (see Figure 1). The latter is especially relevant as it is directly tied to the communication of the true distribution.

9. The variance of the beta distribution is $\sigma^2 = \alpha\beta/(\alpha + \beta)^2(\alpha + \beta + 1)$. The standard deviation σ is more intuitive because it is on the original scale of the variable, in our case vote shares.

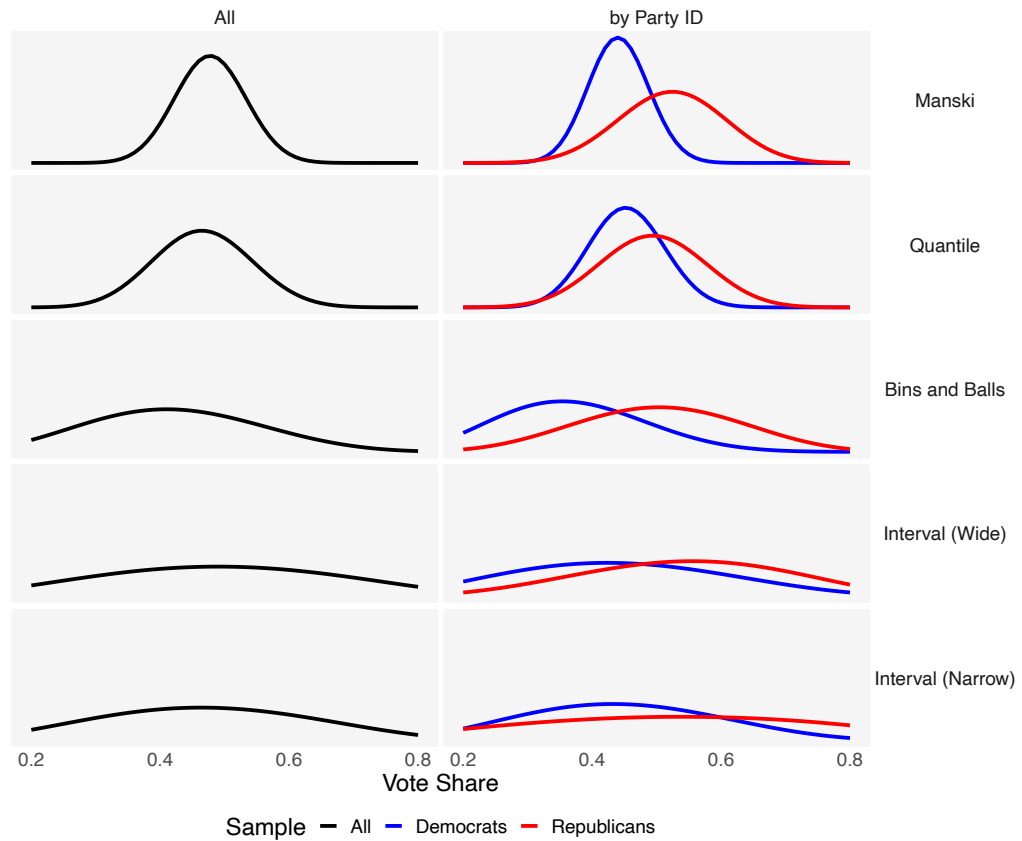


Figure 5. Elicited Beliefs over Trump's Popular Vote Share in November.

what we know about vote shares in US presidential elections, the variances provided by the two Interval questions and Bins and Balls are much too wide to be of any substantive use. The average beliefs elicited by the Interval question have a standard deviation of $\sigma = 17.9$ percent (wide) and $\sigma = 16.3$ percent (narrow) while the beliefs produced by Bins and Balls have standard deviation of $\sigma = 13.1$ percent. In line with the experimental results, both the Manski and Quantile question formats provide reasonable results ($\sigma = 5.4$ and $\sigma = 7.5$ percent, respectively), but the Manski question format again performs best.

The second result is that there are clear partisan differences in the beliefs and expectations about the upcoming 2020 presidential election. We think this is in line with prior literature (Lebo and Cassino 2007; Madson and Hillygus 2019; Kuru, Pasek, and Traugott 2017). Relying on the preferred Manski question format, we find that Republicans have a more optimistic belief about Donald Trump's expected popular vote share (52.3 percent) than Democrats (44.1 percent). At the same time, Republicans are less certain about the election outcome than Democrats. The standard deviation of Republicans' belief is $\sigma = 8.1$ percent compared to only $\sigma = 4.7$ percent for Democrats.

6 Conclusion

This research note has empirically evaluated five different question formats for prior elicitation in the context of online surveys. For each format, we derived the estimators to recover the shape parameters describing respondents beliefs and ran experiments to compare the relative performance of these elicitation methods. We find that a set of questions originally proposed by Manski 2009 performs very well.

This is good news for applied researchers who seek to study citizens' beliefs as distributions. While all five types of elicitation methods are fairly easy to implement, the Manski question is especially straightforward as it only consists of asking people for five numbers (most likely value, lower and upper bound and the probabilities associated with the two bounds). Since it is purely verbal, there is no need for programming - unlike other elicitation methods, such as the Bins and Balls method recently proposed by (Goldstein and Rothschild 2014) which requires programming in JavaScript. In addition, the Manski format seems to perform in a similar fashion across different subgroups defined by political sophistication, which can be a relevant consideration. Finally, there is one caveat that needs mention: we assumed throughout that citizens' beliefs follow a unimodal distribution. While this is a reasonable assumption in many circumstances, it is possible that one would want to elicit multimodal beliefs. In such situation, the Bins and Balls format would allow researchers to do so, but the estimation methods must be adapted accordingly.

Data Availability

Supplementary materials for this article are available on the Political Analysis website. For Data-verse replication materials, see Leemann, Traunmueller, and Stoetzer 2020.

Supplementary Material

(This is dummy text) For supplementary material accompanying this paper, please visit <https://doi.org/10.1017/pan.xxxx.xx>.

References

- Alvarez, R. Michael, and John Brehm. 1997. "Are Americans Ambivalent Towards Racial Policies?" *American Journal of Political Science* 41:345–374.
- Bartels, Larry M. 2002. "Beyond the running tally: Partisan bias in political perceptions." *Political behavior* 24 (2): 117–150.
- Berinsky, Adam J. 2017. "Rumors and health care reform: Experiments in political misinformation." *British Journal of Political Science* 47 (2): 241–262.
- Berinsky, Adam J, Gregory A Huber, and Gabriel S Lenz. 2012. "Evaluating online labor markets for experimental research: Amazon.com's Mechanical Turk." *Political Analysis* 20 (3): 351–368.
- Berinsky, Adam J, Michele F Margolis, and Michael W Sances. 2014. "Separating the shirkers from the workers? Making sure respondents pay attention on self-administered surveys." *American Journal of Political Science* 58 (3): 739–753.
- Bullock, John G. 2009. "Partisan bias and the Bayesian ideal in the study of public opinion." *The Journal of Politics* 71 (3): 1109–1124.
- Coppock, Alexander. 2018. "Generalizing from survey experiments conducted on mechanical Turk: A replication approach." *Political Science Research and Methods*, 1–16.
- Cosmides, Leda, and John Tooby. 1996. "Are Humans Good Intuitive Statisticians after all? Rethinking some Conclusions from the Literature on Judgment under Uncertainty." *cognition* 58 (1): 1–73.
- Garthwaite, Paul H, Joseph B Kadane, and Anthony O'Hagan. 2005. "Statistical methods for eliciting probability distributions." *Journal of the American Statistical Association* 100 (470): 680–701.
- Gerber, Alan, and Donald Green. 1999. "Misperceptions About Perceptual Bias." *Annual Review of Political Science* 2 (1): 189–210.
- Gill, Jeff, and John R Freeman. 2013. "Dynamic elicited priors for updating covert networks." *Network Science* 1 (1): 68–94.

- Gill, Jeff, and Lee D Walker. 2005. "Elicited priors for Bayesian model specifications in political science research." *The Journal of Politics* 67 (3): 841–872.
- Goldstein, Daniel G, and David Rothschild. 2014. "Lay understanding of probability distributions." *Judgment & Decision Making* 9 (1).
- Kuru, Ozan, Josh Pasek, and Michael W Traugott. 2017. "Motivated reasoning in the perceived credibility of public opinion polls." *Public opinion quarterly* 81 (2): 422–446.
- Kynn, Mary. 2008. "The 'heuristics and biases' bias in expert elicitation." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 171 (1): 239–264.
- Lebo, Matthew J, and Daniel Cassino. 2007. "The aggregated consequences of motivated reasoning and the dynamics of partisan presidential approval." *Political Psychology* 28 (6): 719–746.
- Leemann, Lucas, Richard Traunmueller, and Lukas F. Stoetzer. 2020. "Replication Data for: Eliciting Beliefs as Distributions in Online Surveys." *Harvard Dataverse*, <https://doi.org/10.7910/DVN/GEC2LS>.
- Leiter, Andreas Murr, Debra, Ericka Rascón Ramirez, and Mary Stegmaier. 2018. "Social Networks and Citizen Election Forecasting: The More Friends the Better." *International Journal of Forecasting* 34 (2): 235–248.
- Madson, Gabriel J, and D Sunshine Hillygus. 2019. "All the best polls agree with me: Bias in evaluations of political polling." *Political Behavior*, 1–18.
- Manski, Charles F. 2009. *Identification for prediction and decision*. Harvard University Press.
- Mason, Winter, and Siddharth Suri. 2012. "Conducting behavioral research on Amazon's Mechanical Turk." *Behavior research methods* 44 (1): 1–23.
- Morris, David E, Jeremy E Oakley, and John A Crowe. 2014. "A web-based tool for eliciting probability distributions from experts." *Environmental Modelling & Software* 52:1–4.
- Mullinix, Kevin J, Thomas J Leeper, James N Druckman, and Jeremy Freese. 2015. "The generalizability of survey experiments." *Journal of Experimental Political Science* 2 (2): 109–138.
- Murr, Andreas. 2011. "'Wisdom of crowds'? A decentralised election forecasting model that uses citizens' local expectations." *Electoral Studies* 30 (4): 771–783.
- O'Hagan, Anthony, Caitlin E Buck, Alireza Daneshkhah, J Richard Eiser, Paul H Garthwaite, David J Jenkinson, Jeremy E Oakley, and Tim Rakow. 2006. *Uncertain judgements: eliciting experts' probabilities*. John Wiley & Sons.
- Peterson, Cameron, and Alan Miller. 1964. "Mode, median, and mean as optimal strategies." *Journal of Experimental Psychology* 68 (4): 363.
- Savage, Leonard J. 1971. "Elicitation of personal probabilities and expectations." *Journal of the American Statistical Association* 66 (336): 783–801.
- Schlag, Karl H, James Tremewan, and Joël J Van der Weele. 2015. "A penny for your thoughts: a survey of methods for eliciting beliefs." *Experimental Economics* 18 (3): 457–490.
- Spetzler, Carl S, and Carl-Axel S Stael von Holstein. 1975. "Exceptional paper—Probability encoding in decision analysis." *Management science* 22 (3): 340–358.
- Thomas, Kyle A, and Scott Clifford. 2017. "Validity and mechanical turk: An assessment of exclusion methods and interactive experiments." *Computers in Human Behavior* 77:184–197.
- Tversky, Amos, and Daniel Kahneman. 1971. "Belief in the Law of Small Numbers." *Psychological Bulletin* 76 (2): 105.
- Tversky, Amos, and Daniel Kahneman. 1973. "Availability: A Heuristic for Judging Frequency and Probability." *Cognitive Psychology* 5 (2): 207–232.

———. 1974. “Judgment under uncertainty: Heuristics and biases.” *science* 185 (4157): 1124–1131.

Wallsten, Thomas S, and David V Budescu. 1983. “State of the Art – Encoding Subjective Probabilities: A Psychological and Psychometric Review.” *Management Science* 29 (2): 151–173.

Zaller, John, and Stanley Feldman. 1992. “A Simple Theory of the Survey Response: Answering Questions and Revealing Preferences.” *American Journal of Political Science* 36:579–616.

APPENDIX

Lucas Leemann[†] Lukas F. Stoetzer[‡] Richard Traunmüller[§]

August 21, 2020

Contents

A1 Question Formats	2
A2 Estimation	4
A2.1 Likelihood for the Quantile Question	4
A2.2 Likelihood for the Interval Question	5
A2.3 Likelihood for the Manski Question	6
A2.4 Likelihood for the Bins and Balls Question	6
A3 Balance	8
A4 Evaluating Adequacy Check for the Quantile Question	10
A5 No Screening	11
A5.1 Individual Beliefs	13
A5.2 Results	14
A6 Sub-Samples Political Interest	16
A7 Additional Survey Trump Vote Share	18
A8 Timing of Elicitation Methods	19

* This is the document intended as online appendix for the manuscript *Eliciting Beliefs as Distributions in Online Surveys*

[†]Department of Political Science, University of Zurich, Switzerland. Email: leemann@ipz.uzh.ch, URL: <http://www.lucasleemann.ch>

[‡]Department of Political Science, Humboldt University of Berlin, Germany. Email: lukas.stoetzer@hu-berlin.de, URL: <http://lukas-stoetzer.org/>

[§]University of Mannheim, Germany. URL: www.richardtraunmueller.com, Email: traunmueller@uni-mannheim.de

A1 Question Formats

• Quantile Question

1. We have just shown you election results for similar districts. Now, we want to know what your expectations are.
Can you determine the median? This is the value where the vote share of party A is equally likely to be less than or larger than this value.
2. Imagine you were told that the actual result was below your median value.
Can you determine a new value, such that the vote share of party A is equally likely to be *between 0 percent and the new value* or *between the new value and the median value*?
3. Imagine you were told that the actual result was above your median value.
Can you determine yet another value, such that the vote share of party A is equally likely to be *between the median and this new value* or *between this new value and 100 percent*?
4. At the end, respondents are shown the four ranges (0-25th, 25th-50th, 50th-75th, 75th-100) and asked whether a random draw is equally likely to occur in each of them. If not, respondents can go and adjust their responses.
 - (a) Consider the following four intervals: $[0, P_{25}]$, $[P_{25}, P_{50}]$, $[P_{50}, P_{75}]$, $[P_{75}, 100]$.
Is it equally likely that party A's vote share will fall in any of these intervals?
(P_{XY} indicates the respondent's XY percentile.)

• Interval Question (Narrow and Wide) This question comes in two versions – a wide and a narrow version.

1. We have just shown you election results for similar districts. Now, we want to know what your expectations for party A's vote share are.
What is the most likely vote share of party A? Please give your response in percentage points.
2. What is the probability that party A will receive a vote share of less than 40 percent? (*45% in narrow format*)
3. What is the probability that party A will receive a vote share of more than 60 percent? (*55% in narrow format*)

• Manski Question

1. We have just shown you election results for similar districts. Now, we want to know what your expectations for party A's vote share are.
What is the most likely vote share of party A? Please give your response in percentage points.
2. What do you think is a likely range of the vote share that party A will receive?
Please indicate the lower bound in percentage points.
3. Now, please indicate the upper bound in percentage points.

4. What is the probability that party A will get a vote share of less than (*lower value indicated by R*) percent?
5. What is the probability that party A will get a vote share of more than (*upper value indicated by R*) percent?

- **Bins and Balls**

We implemented this question by inserting Java script into the survey software - a step that should be easily replicated by anybody. We also provide the JS code here (link) and have annotated it so that it can be adapted easily.

Figure A1 shows the full question as it is presented to respondents. By clicking on + and − buttons, the various bins can be filled or emptied. Each respondent allocates 100 balls into these bins.

Imagine 100 elections that are being held in a given district. What are the likely vote shares for party A over these 100 elections? Please click on the buttons below to describe the distribution of vote shares you would expect.

There are 100 left.

− +	25-29%	0
− +	30-34%	0
− +	35-39%	0
− +	40-44%	0
− +	45-49%	0
− +	50-54%	0
− +	55-59%	0
− +	60-64%	0
− +	65-69%	0
− +	70-74%	0
− +	75-79%	0
− +	80-84%	0

Figure A1: Screenshot of Balls and Bins

A2 Estimation

To estimate beliefs from the different question formats, we develop statistical models that permit us to estimate the parameters of respondents' belief distributions. In the following, we describe the Likelihoods that model the observed outcomes given parametric belief distributions for the different question formats.

A2.1 Likelihood for the Quantile Question

For the quantile question, we observe three outcomes for each respondent: The lower quartile, the median value, and the upper quartile of the respondent's belief. We denote them with $y_i = [y_{i1}, y_{i2}, y_{i3}]$, where $i \in (1, \dots, N)$ are respondents and k refers to the different quartile questions $k \in (1, 2, 3)$. We assume that these values are observed with measurement error, such that:¹

$$y_{ik} \sim \mathcal{N}(\mu_{ik}, \sigma^2). \quad (1)$$

To estimate a respondent's average belief, we assume a parametric distribution and estimate the parameters of the distribution to closely mimic the expected observed indicators. To map the beliefs *of/about?* the measured values, we require the quantile function of the belief distribution. Because the sampling space of our experiment is bound between 0 and 1, we employ a Beta distribution as our parametric distribution. We denote $Q^{-1}(q_k, \alpha, \beta)$ as the quantile function of the beta distribution. The two shape parameters α and β define the expectation and the variance of the belief. The distribution is then linked to the expectation of the observed values. If we denote the three quartiles with $q_k = [0.25, 0.5, 0.75]$, we can write:

$$\mu_k = Q^{-1}(q_k, \alpha, \beta) \quad (2)$$

With this model, we define the likelihood of a respondent's observed answers as:

$$L(\alpha, \beta, \sigma^2 \mid y_i) = \prod_{k=1}^3 \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[\frac{-(y_{ik} - \mu_k)^2}{2\sigma^2} \right] \quad (3)$$

Maximizing the Likelihood for each individual would involve minimizing the squared distance between the observed quartile measurements and the shape parameters of the beta-distribution that generate the expected quartiles. If we assume that individual responses are identical and independently distributed, we can further write the Likelihood for the full sample as:

$$L(\alpha, \beta, \sigma^2 \mid \mathbf{Y}) = \prod_{i=1}^N \prod_{k=1}^3 \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[\frac{-(y_{ik} - Q^{-1}(q_k, \alpha, \beta))^2}{2\sigma^2} \right], \quad (4)$$

¹We assume that the measurement errors are normally distributed with the same error variance and no covariance between the errors.

where \mathbf{Y} is a $(N \times K)$ matrix with all respondents' responses $[y_1, \dots, y_N]'$. The function is maximized with respect to the parameters α, β, σ using R's `optim` function.

A2.2 Likelihood for the Interval Question

We observe three values for the interval question. Respondents report the mean value of their beliefs and the probabilities of observing a value below and above a certain threshold. We denote the mean with y_i and the two ($k \in (1, 2)$) probabilities with p_{i1}, p_{i2} . The interval values depend on the question format and are denoted with $c = [c_1, c_2]$, where in the wide version $c = [40\%, 60\%]$ and in the narrow version $c = [45\%, 55\%]$. We assume that the values are measured with normal measurement error.

$$y_i \sim \mathcal{N}(\mu_y, \sigma_y^2) \quad (5)$$

$$p_{i1} \sim \mathcal{N}(\mu_{p1}, \sigma_p^2) \quad (6)$$

$$p_{i2} \sim \mathcal{N}(\mu_{p2}, \sigma_p^2) \quad (7)$$

The expectations μ_y are calculated from the assumed parametric belief distribution. Here, we use the same distribution as in the data-generating process - a beta distribution. The beta distribution is relatively flexible and well-suited for our example with vote shares being constrained on the unit interval. It is generally possible to use other parametric distributions, like a normal distribution, instead. In practical applications, it would be sensible to try different distributions and compare their relative fit. The expectation for the mean from the beta are given by the two shape parameters α and β :

$$\mu_y = \frac{\alpha}{\alpha + \beta} \quad (8)$$

The expected probabilities are given by the CDF of the beta distribution, which we denote with $Q(\cdot, \alpha, \beta)$.

$$\mu_{p1} = Q(c_1, \alpha, \beta) \quad (9)$$

$$\mu_{p2} = 1 - Q(c_2, \alpha, \beta) \quad (10)$$

With this model, we can define the Likelihood for the observed data of N respondents $Y = [[y_1, p_{i1}, p_{i2}]', \dots, [y_N, p_{N1}, p_{N2}]']'$. We assume that all responses are identically and independently distributed, which yields the following Likelihood:

$$L(\alpha, \beta, \sigma_y^2, \sigma_p^2 \mid \mathbf{Y}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp \left[-\frac{(y_i - \mu_y)^2}{2\sigma_y^2} \right] \prod_{k=1}^2 \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp \left[-\frac{(p_k - \mu_{p_k})^2}{2\sigma_p^2} \right]. \quad (11)$$

To obtain MLE estimates of the parameters, the function is also maximized using R's `optim` function. The obtained estimates yield an estimate of the average beliefs under a specific condition. The goal is to identify the question format that will yield estimates that come closest to the true values.

A2.3 Likelihood for the Manski Question

For the Manski Question, we observe five measures of respondents' beliefs. We measure three $k \in 1, 2, 3$ values: the mean value (which we denote with y_{i1}), and the lower and the upper bound values (which we denote with y_{i2} and y_{i3} , respectively). In addition, we measure two probabilities of observing values below and above the bounds (which we denote with p_{i1} and p_{i2}). We assume that the values are measured with error and that the errors are identical and independently normally distributed.

$$y_{ik} \sim \mathcal{N}(\mu_k, \sigma^2) \quad (12)$$

the expectations μ_k are calculated from the assumed parametric distribution of respondents' beliefs. In our analysis, we work with the beta distribution, which yields a simple expression for the mean value. Given the assigned probabilities, the observed lower and upper bounds can be calculated from the quantile function of the Beta distribution, which we denote as $Q^{-1}(\cdot, \alpha, \beta)$. The expectations of the measurement model are given by:

$$\mu_{i1} = \frac{\alpha}{\alpha + \beta} \quad (13)$$

$$\mu_{i2} = Q^{-1}(p_{i1}, \alpha, \beta) \quad (14)$$

$$\mu_{i3} = 1 - Q^{-1}(p_{i2}, \alpha, \beta) \quad (15)$$

The Likelihood is given by the normal measurement error and the respective expectation-generating functions. We collapse the measured values and the probabilities in a matrix ($\mathbf{Y} = [[y_{11}, y_{12}, y_{13}, p_{11}, p_{12}]', \dots, [y_{N1}, y_{N2}, y_{N3}, p_{N1}, p_{N2}]']'$). Assuming that the observed values are independent allows us to write the Likelihood as:

$$L(\alpha, \beta, \sigma \mid \mathbf{Y}) = \prod_{i=1}^N \prod_{k=1}^3 \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_{ik} - \mu_{ik})^2}{2\sigma^2} \right], \quad (16)$$

which is numerically maximized with respect to the parameters using R's optim function.

A2.4 Likelihood for the Bins and Balls Question

The bins and balls question has a slightly different structure compared to the quantile questions. For this question, we observe the number of balls a respondent decides to place into K bins that each covers an exclusive interval. The intervals are given by ordered cut points c_1, \dots, c_C . There is one cut-point more than categories $C = K + 1$, as the question format can have lower and upper bounds.² The number of balls out of $B = 100$ that a respondent places in a bin is denoted with y_{ik} . We assume that the measured placements are binomially distributed, with a certain probability π_k .

² C is the number of cut-points which, in our question format, is $C = 13$, and the corresponding cut points are 0.25, 0.3, 0.35, 0.4, 0.45, \dots , 0.85.

$$y_{ik} \sim \mathcal{B}(\pi_k, B) \quad (17)$$

The probabilities are calculated from the CDF of the assumed parametric belief distribution. The CDF of the Beta distribution is given by $Q(\cdot, \alpha, \beta)$. With this, we calculate the probability that a respondent places balls in each bin, as:

$$\pi_k = Q(c_{k+1}, \alpha, \beta) - Q(c_k, \alpha, \beta). \quad (18)$$

Assuming that the observed values are conditionally independent, combining all observed placements $\mathbf{Y} = [[y_{11}, \dots, y_{1K}]', \dots, [y_{N1}, \dots, y_{NK}]']'$ yields the following Likelihood:

$$L(\alpha, \beta \mid \mathbf{Y}) = \prod_{i=1}^N \prod_{k=1}^K \binom{B}{y_{ik}} \pi_k^{y_{ik}} (1 - \pi_k)^{B-y_{ik}} \quad (19)$$

We numerically maximize the likelihood of obtaining MLE estimates of the shape parameters.

A3 Balance

The following four tables present balance checks in terms of covariate averages and standard deviations for each experimental condition. These checks are based on the full data before reducing the data set only to observations which passed both attention checks. Please note that we refrain from the ill-advised practice of statistically *testing* for mean differences (Mutz, 2011). Overall, we find treatment conditions to be well balanced.

Format	n	Female	Age	University	Political Interest
Quantile	196	0.45 (0.50)	41.40 (11.73)	0.61 (0.49)	3.03 (0.76)
Interval (Wide)	205	0.44 (0.50)	40.33 (13.76)	0.60 (0.49)	3.00 (0.82)
Interval (Narrow)	205	0.48 (0.50)	41.39 (12.26)	0.57 (0.50)	3.09 (0.78)
Manski	201	0.49 (0.50)	41.69 (11.64)	0.60 (0.49)	3.04 (0.81)
Bins and Balls	189	0.49 (0.50)	40.81 (11.63)	0.61 (0.49)	2.99 (0.81)

Table A1: Balance Check for a Symmetric Distribution. Means and Standard Deviations in Parentheses.

Format	n	Female	Age	University	Political Interest
Quantile	100	0.43 (0.50)	38.04 (11.74)	0.67 (0.47)	2.93 (0.84)
Interval (Wide)	102	0.40 (0.49)	39.29 (12.29)	0.53 (0.50)	2.92 (0.86)
Interval (Narrow)	97	0.40 (0.49)	40.66 (12.66)	0.60 (0.49)	2.95 (0.85)
Manski	106	0.42 (0.50)	41.20 (12.39)	0.56 (0.50)	2.93 (0.80)
Bins and Balls	102	0.41 (0.49)	41.35 (13.58)	0.65 (0.48)	3.08 (0.83)

Table A2: Balance Check for a Symmetric Distribution with a Large Variance. Means and Standard Deviations in Parentheses

Format	n	Female	Age	University	Political Interest
Quantile	201	0.48 (0.50)	41.29 (12.24)	0.63 (0.48)	3.00 (0.73)
Interval (Wide)	197	0.40 (0.49)	40.43 (11.61)	0.59 (0.49)	3.04 (0.77)
Interval (Narrow)	206	0.47 (0.50)	39.47 (12.13)	0.59 (0.49)	3.02 (0.80)
Manski	203	0.46 (0.50)	39.93 (12.30)	0.59 (0.49)	3.04 (0.78)
Bins and Balls	196	0.43 (0.50)	41.44 (12.07)	0.64 (0.48)	3.10 (0.80)

Table A3: Balance Check for an Asymmetric Distribution. Means and Standard Deviations in Parentheses.

Format	n	Female	Age	University	Political Interest
Quantile	102.00	0.54 (0.50)	41.49 (12.40)	0.72 (0.45)	3.09 (0.76)
Interval (Wide)	104.00	0.44 (0.50)	38.80 (10.72)	0.72 (0.45)	2.88 (0.75)
Interval (Narrow)	101.00	0.45 (0.50)	39.73 (12.18)	0.65 (0.48)	2.99 (0.83)
Manski	98.00	0.48 (0.50)	41.00 (13.26)	0.68 (0.47)	2.91 (0.90)
Bins and Balls	95.00	0.53 (0.50)	42.15 (12.26)	0.49 (0.50)	3.00 (0.77)

Table A4: Balance Check for an Asymmetric Distribution with Large Variance. Means and Standard Deviations in Parentheses.

A4 Evaluating Adequacy Check for the Quantile Question

Variance	Distribution	AdequacyCheck	alpha	beta	KL	N
Large Variance	Asymmetric	Yes	49.08	27.64	0.13	118
Large Variance	Asymmetric	No	158.50	91.07	0.98	130
Large Variance	Symmetric	Yes	48.41	47.90	0.01	119
Large Variance	Symmetric	No	34.21	34.59	0.07	115

Table A5: Estimates for the Elicited Beliefs. Comparing Quantile with and without Adequacy Check for the Large Variance Scenarios

A5 No Screening

These two tables present the same information that is shown in Table 2. The data here is based on all results, i.e the raw data before reducing the data set only to observations which pass both attention checks.

method	alpha	beta	KL	lr	N
Quantile	45.54	47.29	0.07	0.14	100
Bins and Balls	15.27	15.04	0.10	0.00	102
Manski	13.57	13.34	0.13	0.00	106
Interval (Wide)	8.80	8.53	0.27	0.00	97
Interval (Narrow)	7.13	6.78	0.36	0.00	102

(a) Symmetric, Large Variance

method	alpha	beta	KL	lr	N
Quantile	59.69	60.15	0.00	0.85	196
Manski	50.73	49.12	0.02	0.04	201
Bins and Balls	26.76	26.82	0.13	0.00	189
Interval (Wide)	15.98	15.40	0.31	0.00	205
Interval (Narrow)	11.56	11.36	0.43	0.00	205

(c) Symmetric, Small Variance

method	alpha	beta	KL	lr	N
Manski	22.79	13.46	0.12	0.00	98
Bins and Balls	11.49	7.16	0.24	0.00	95
Interval (Wide)	5.82	3.34	0.42	0.00	101
Quantile	60.58	35.75	0.54	0.00	102
Interval (Narrow)	3.40	1.82	0.67	0.00	104

(b) Asymmetric, Large Variance

method	alpha	beta	KL	lr	N
Manski	38.79	20.29	0.05	0.00	203
Quantile	56.95	31.47	0.10	0.00	201
Bins and Balls	18.27	12.45	0.54	0.00	196
Interval (Wide)	7.85	3.98	0.60	0.00	206
Interval (Narrow)	6.56	3.12	0.70	0.00	197

(d) Asymmetric, Small Variance

Table A6: Estimates for Elicited Beliefs. No Screen

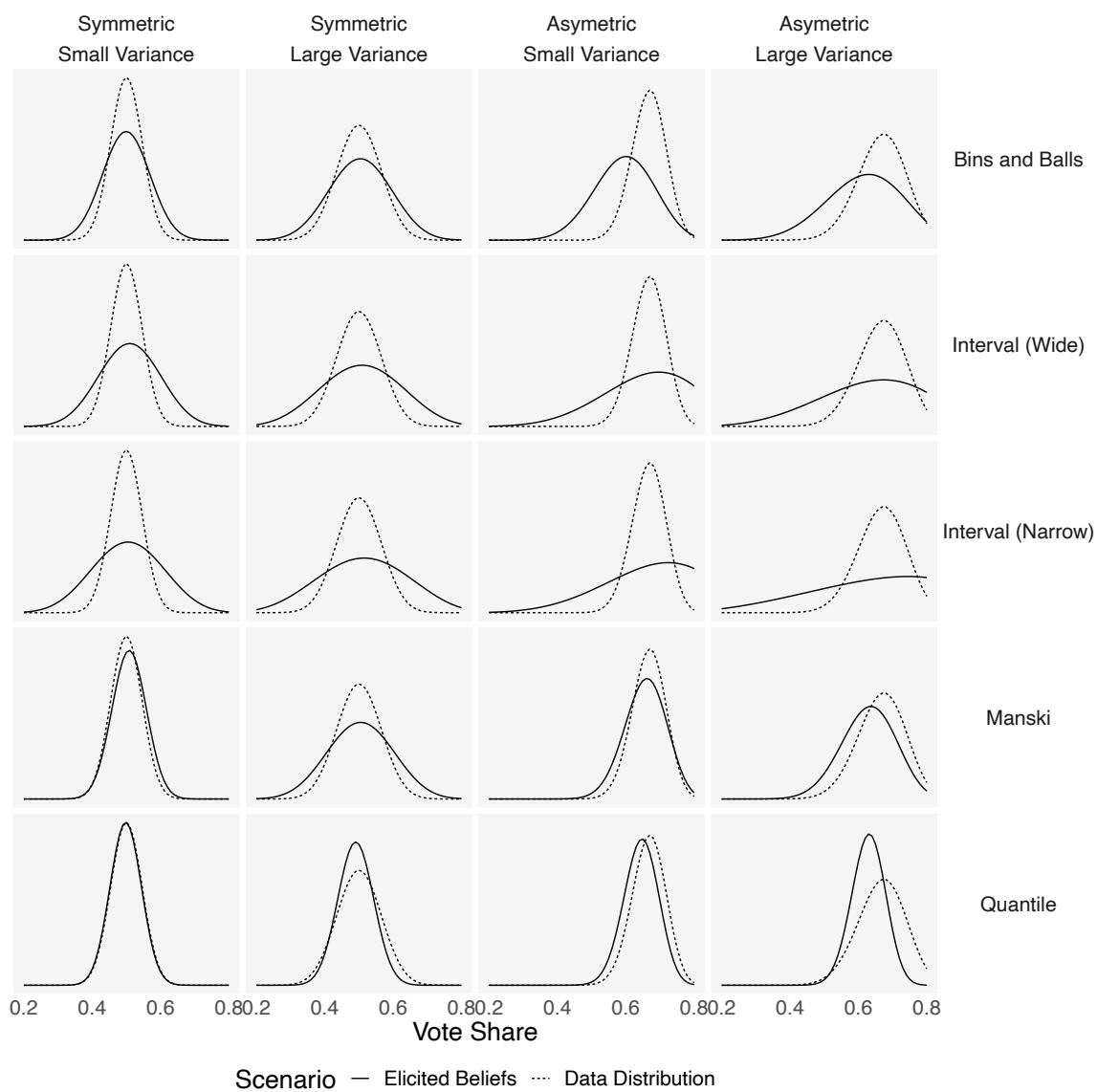


Figure A2: Comparison of Question Formats. No screen. The dotted line indicates the true distribution and the black solid line shows the average of the elicited distributions.

A5.1 Individual Beliefs

The question formats and estimation method can also be used to obtain individual beliefs. We use the same Maximum Likelihood approach as described in section A2, but allow for individual shape parameters: $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]$ and $\boldsymbol{\beta} = [\beta_1, \dots, \beta_N]$. To illustrate, consider the Individual Likelihood for the Quantile question:

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}, \sigma \mid \mathbf{Y}) = \prod_{i=1}^N \prod_{k=1}^3 \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_{ik} - Q^{-1}(q_k, \alpha_i, \beta_i))^2}{2\sigma^2} \right], \quad (20)$$

where we now introduce a subscript for the shape parameters of the Beta quantile function $Q^{-1}(q_k, \alpha_i, \beta_i)$.

We obtain estimates for respondent-specific shape parameters by numerically maximizing the Likelihood function. We first estimate the shape parameters for each respondent, and afterwards estimate the error variance terms for the Likelihood function. We repeat until convergence in the error variances.³

³Some response patterns do not yield estimates of sensible shape parameters. For example, if a respondent reports a lower quartile of 0.50 and a Median of 0.45, the maximization of the function will be impossible. We exclude respondents with such inconsistent answering behaviors.

A5.2 Results

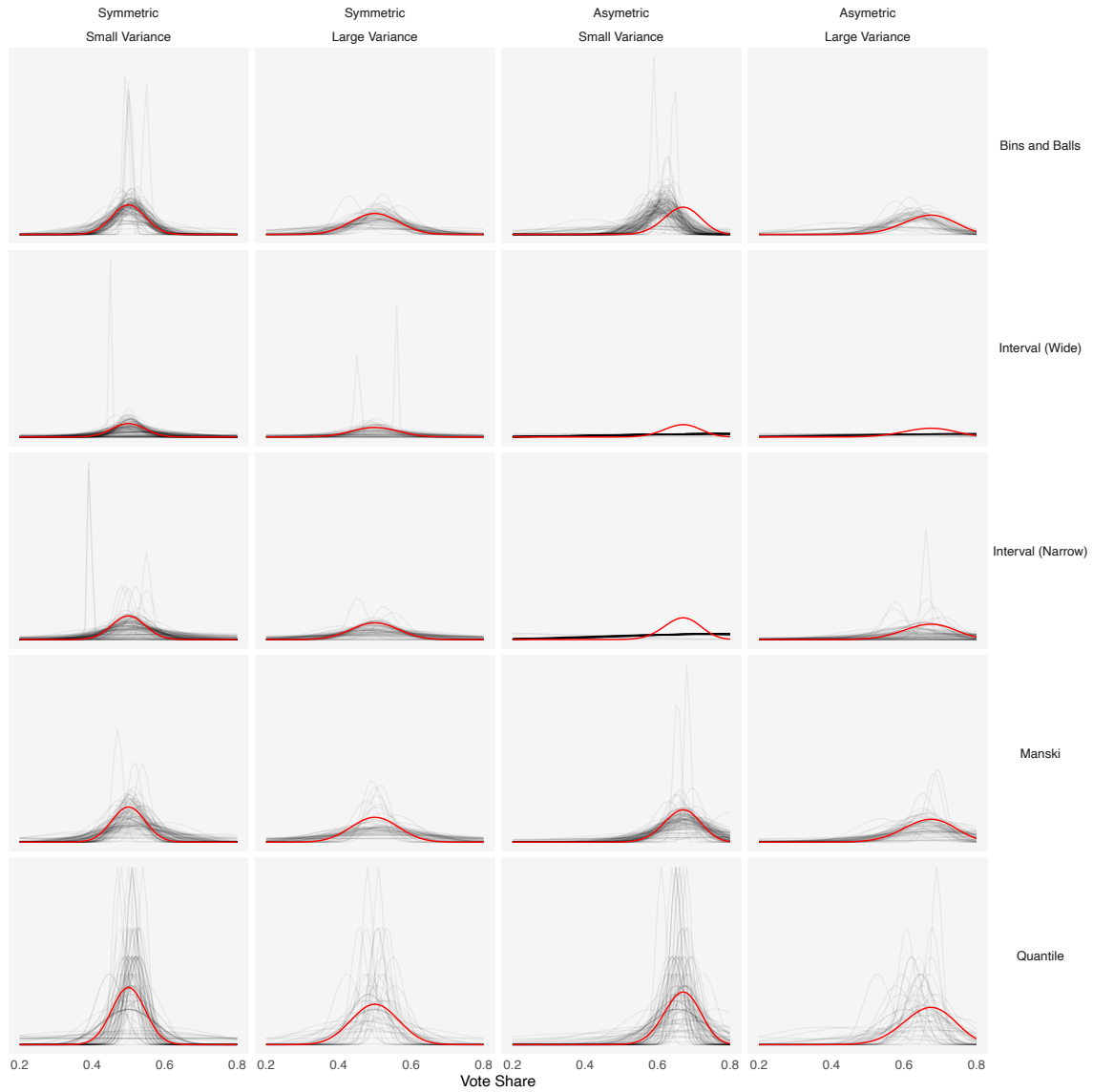


Figure A3: Individual Beliefs. The grey lines indicate individual elicited beliefs. The red line indicates the true distribution.

type	method	median	qlow	qhigh
Symmetric	Bins and Balls	0.13	0.06	0.26
Symmetric	Interval (Narrow)	0.27	0.05	0.81
Symmetric	Manski	0.36	0.23	0.66
Symmetric	Interval (Wide)	0.45	0.14	1.33
Symmetric	Quantile	0.50	0.15	1.92

(a) Symmetric, Large Variance

type	method	median	qlow	qhigh
Symmetric	Bins and Balls	0.10	0.03	0.23
Symmetric	Manski	0.25	0.11	0.45
Symmetric	Interval (Wide)	0.36	0.12	1.27
Symmetric	Interval (Narrow)	0.37	0.11	1.11
Symmetric	Quantile	0.48	0.17	1.17

(c) Symmetric, Small Variance

type	method	median	qlow	qhigh
Asymetric	Bins and Balls	0.27	0.12	0.62
Asymetric	Manski	0.27	0.09	0.53
Asymetric	Quantile	0.51	0.16	1.48
Asymetric	Interval (Wide)	0.64	0.63	0.73
Asymetric	Interval (Narrow)	0.84	0.37	2.26

(b) Asymmetric, Large Variance

type	method	median	qlow	qhigh
Asymetric	Manski	0.15	0.05	0.38
Asymetric	Quantile	0.44	0.14	0.91
Asymetric	Interval (Wide)	0.96	0.95	0.96
Asymetric	Interval (Narrow)	0.96	0.95	0.96
Asymetric	Bins and Balls	0.99	0.55	1.54

(d) Asymmetric, Small Variance

Table A7: KL Divergence for Individual Beliefs in different Scenarios

method	median	qlow	qhigh
Manski	0.25	0.09	0.49
Bins and Balls	0.30	0.09	0.90
Quantile	0.48	0.14	1.20
Interval (Wide)	0.92	0.43	0.96
Interval (Narrow)	0.93	0.25	0.98

Table A8: KL Divergence for Individual Beliefs over different Scenarios

A6 Sub-Samples Political Interest

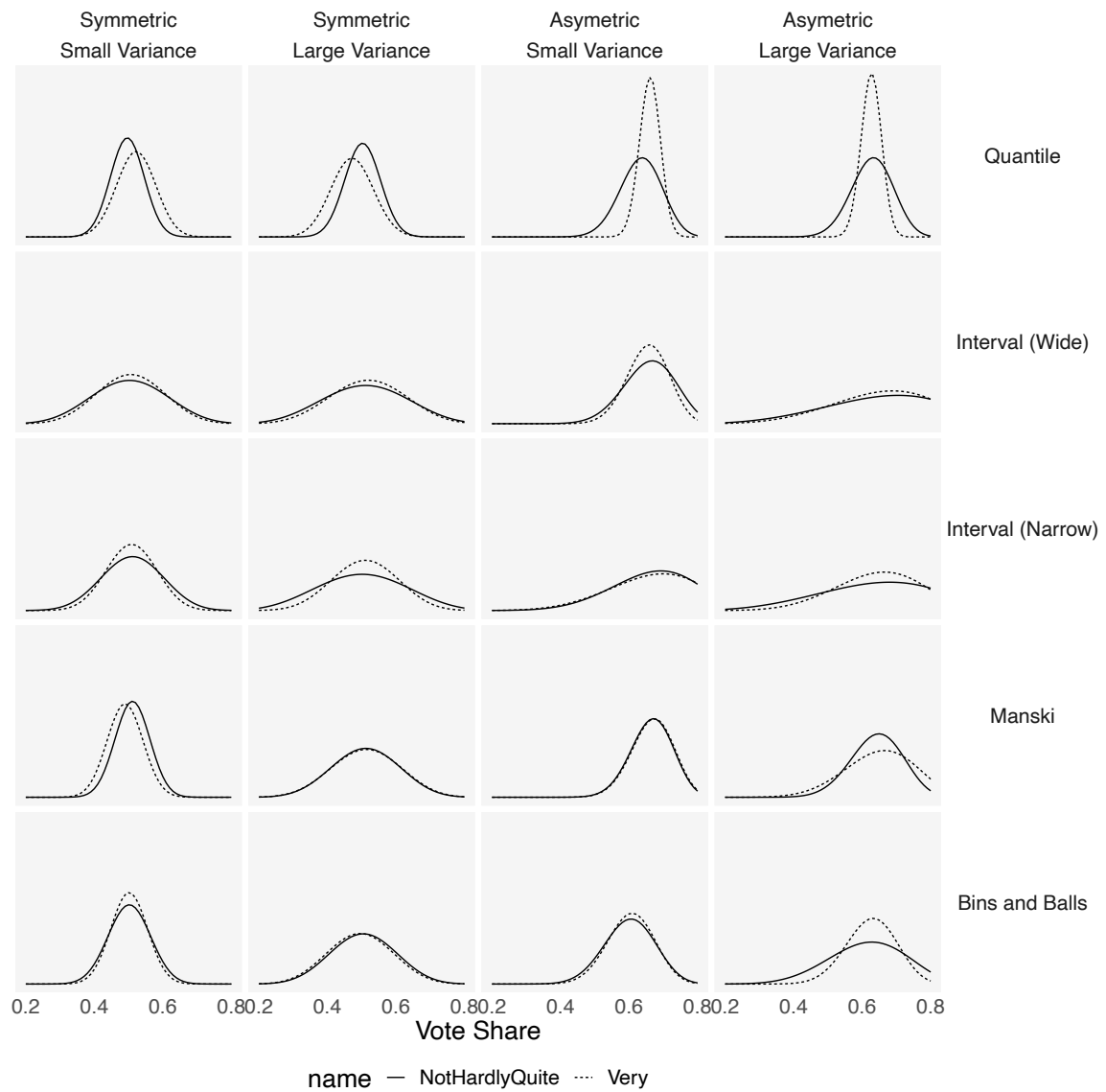


Figure A4: Estimated Beliefs for Sub-Groups of Political Interest

method	NotHardlyQuite_alpha	NotHardlyQuite_beta	NotHardlyQuite_KL	Very_alpha	Very_beta	Very_KL
Quantile	46.80	46.33	0.06	31.01	34.63	0.10
Interval (Wide)	8.11	7.79	0.31	10.56	9.83	0.23
Interval (Narrow)	7.16	7.12	0.35	13.78	13.24	0.13
Manski	13.13	12.55	0.15	12.62	11.98	0.16
Bins and Balls	13.65	13.40	0.13	13.57	13.81	0.12

(a) Symmetric, Large Variance

method	NotHardlyQuite_alpha	NotHardlyQuite_beta	NotHardlyQuite_KL	Very_alpha	Very_beta	Very_KL
Quantile	39.09	23.07	0.24	166.57	99.32	2.70
Interval (Wide)	5.01	2.68	0.49	6.75	3.60	0.36
Interval (Narrow)	5.04	2.90	0.48	9.27	5.17	0.25
Manski	25.16	14.05	0.06	13.69	7.38	0.12
Bins and Balls	10.98	6.91	0.26	26.79	16.02	0.16

(b) Asymmetric, Large Variance

method	NotHardlyQuite_alpha	NotHardlyQuite_beta	NotHardlyQuite_KL	Very_alpha	Very_beta	Very_KL
Quantile	51.63	52.38	0.01	40.24	36.87	0.12
Interval (Wide)	10.18	10.08	0.48	13.15	12.84	0.38
Interval (Narrow)	15.88	15.20	0.32	23.63	22.84	0.18
Manski	49.91	47.72	0.03	45.31	47.19	0.03
Bins and Balls	33.50	33.21	0.07	44.41	44.11	0.02

(c) Symmetric, Small Variance

method	NotHardlyQuite_alpha	NotHardlyQuite_beta	NotHardlyQuite_KL	Very_alpha	Very_beta	Very_KL
Quantile	39.07	22.54	0.17	159.98	82.26	0.38
Interval (Wide)	24.78	12.82	0.15	38.88	20.56	0.06
Interval (Narrow)	9.79	4.90	0.50	8.38	4.17	0.57
Manski	38.64	19.43	0.04	38.58	19.17	0.04
Bins and Balls	25.90	17.20	0.50	30.70	20.06	0.48

(d) Asymmetric, Small Variance

method	NotHardlyQuite_KL	Very_KL
Quantile	0.12	0.83
Interval (Wide)	0.36	0.26
Interval (Narrow)	0.41	0.28
Manski	0.07	0.09
Bins and Balls	0.24	0.20

(e) Summary Kullback-Leibler divergence over four applications

Table A9: Estimates for Sub-groups of Political Interest

A7 Additional Survey Trump Vote Share

In section 5 of the paper, we show the results from an additional survey where we ask respondents about their beliefs regarding Donald Trump’s popular vote share in November.

[Table A10](#) shows the beliefs according to the different formats. In addition, the results are also shown for sub-samples of Republicans and Democrats.

Table A10: Parameter Estimates & Moments of Belief Distribution

	α	β	mean	variance
Manski	39.70	43.38	0.48	0.00
Manski (D)	49.47	62.74	0.44	0.00
Manski (R)	19.34	17.61	0.52	0.01
Quantile	19.91	22.81	0.47	0.01
Quantile (D)	32.39	39.05	0.45	0.00
Quantile (R)	18.57	19.00	0.49	0.01
Bins & Balls	5.57	7.58	0.42	0.02
Bins & Balls (D)	6.43	10.93	0.37	0.01
Bins & Balls (R)	7.46	7.35	0.50	0.02
Interval wide	3.36	3.46	0.49	0.03
Interval wide (D)	3.70	4.67	0.44	0.03
Interval wide (R)	5.03	4.22	0.54	0.02
Interval narrow	3.94	4.40	0.47	0.03
Interval narrow (D)	4.47	5.56	0.45	0.02
Interval narrow (R)	2.39	2.19	0.52	0.04

A8 Timing of Elicitation Methods

The time variable accounts for time for the *entire* survey. As all respondents have the same introduction questions and identical demographic questions, the remaining differences are due to the different belief elicitation question formats. Since these additional variables are constant in all elicitation methods, the timing variable gives us a sense of the relative performance of the five elicitation methods. However, a simple F-test reveals that these differences are not statistically significant.

Table A11: Median amount of time spent per elicitation in seconds

Elicitation Method	Experiment	Trump Vote
Manski Question	170.0	132.0
Quantile Question	200.0	162.0
Interval Question Wide	149.0	106.0
Interval Question Narrow	157.0	108.5
Bins and Balls	199.0	155.0
F-value	0.89	1.56
Pr(>F)	0.4682	0.1833

References

Mutz, Diana C. 2011. *Population-Based Survey Experiments*. Princeton University Press.